



Introduction to Bioinformatics

3. DNA editing and contig assembly

Benjamin F. Matthews

United States Department of Agriculture
Soybean Genomics and Improvement
Laboratory

Beltsville, MD 20708

matthewb@ba.ars.usda.gov



What we will cover today

- DNA editing
 - Phred
- Sequence assembly (Contig building)
 - Phrap
 - Consed
 - CAP3
 - DNA Star - commercial software
 - <http://www.phrap.org/>



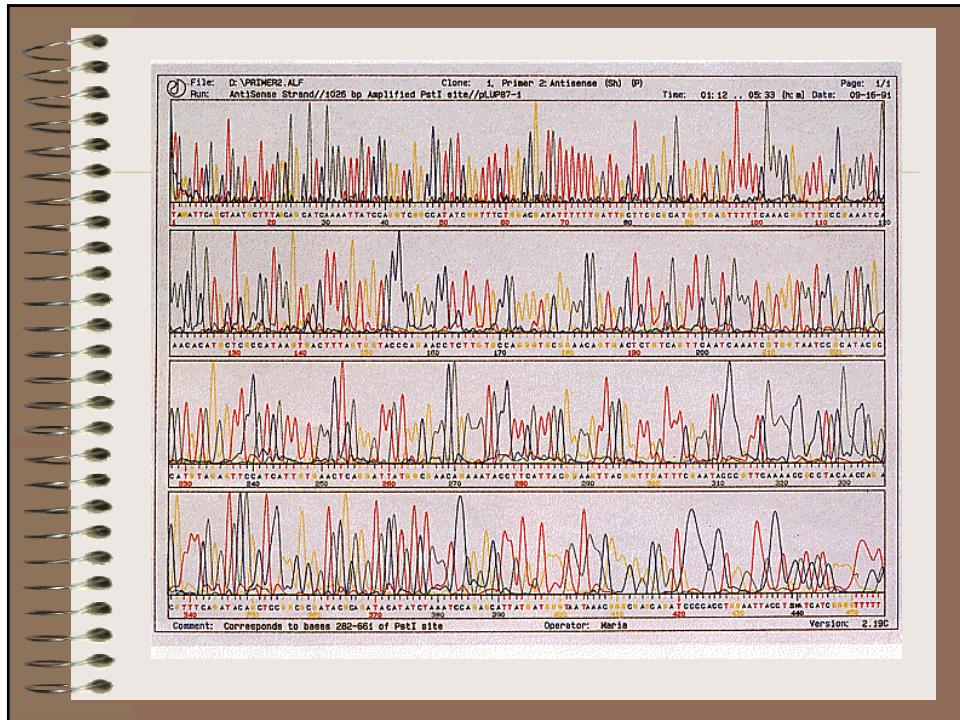
What we will cover today

- DNA Sequencing software
- DNA sequence assembly
- Similarity searching with a DNA sequence
- BLAST



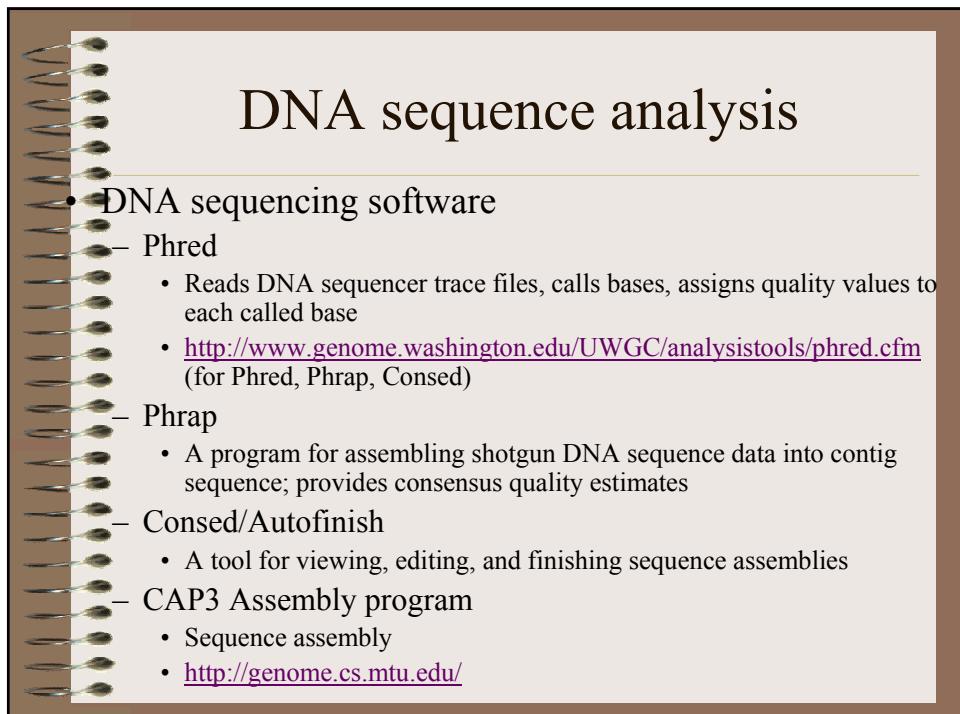
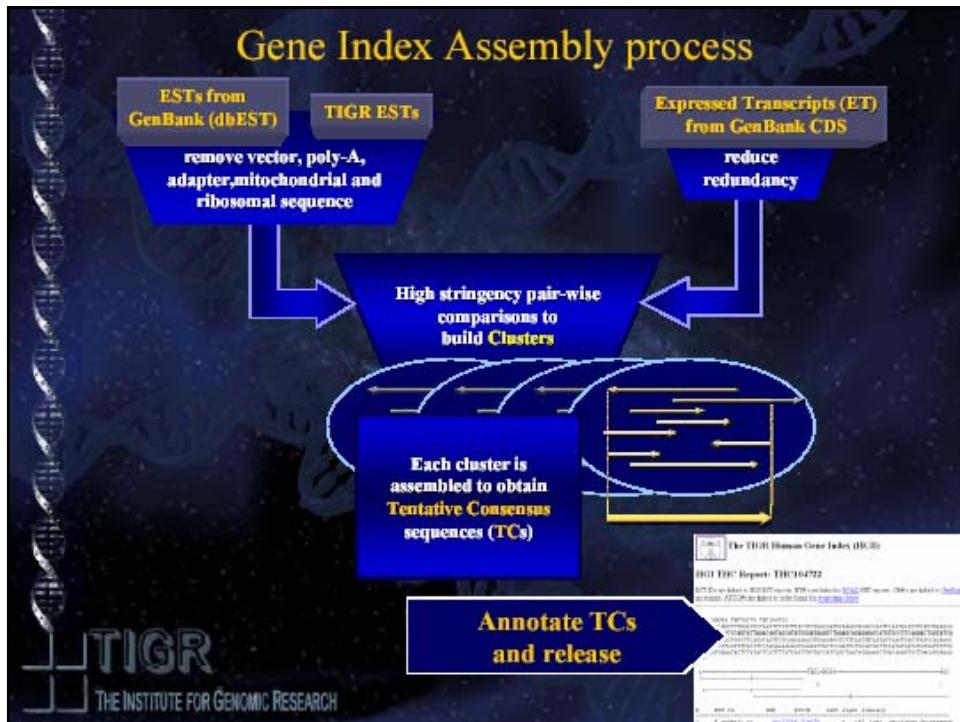
You cloned a cDNA

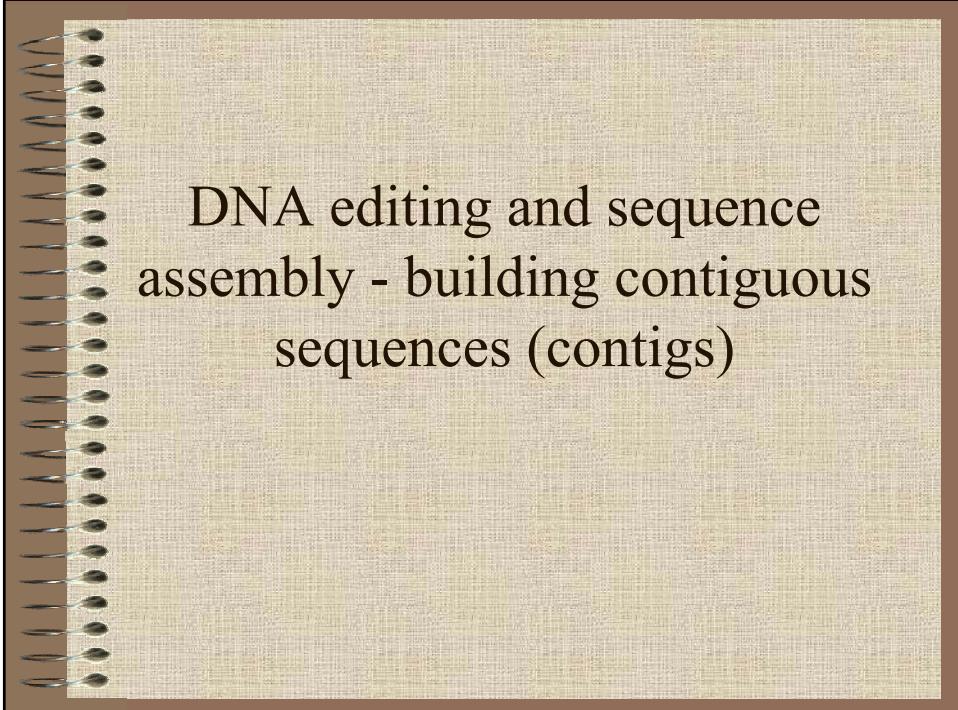
- Isolated mRNA
- Reverse transcribed
- Placed into vector
- Transformed and grew bacteria
- Harvested plasmid
- Sequenced insert



DNA sequence analysis

- Is full-length cDNA cloned?
- What are its properties
- What is function of encoded protein?
- Are there family members?
- Is it cloned from other organisms?





DNA editing and sequence assembly - building contiguous sequences (contigs)

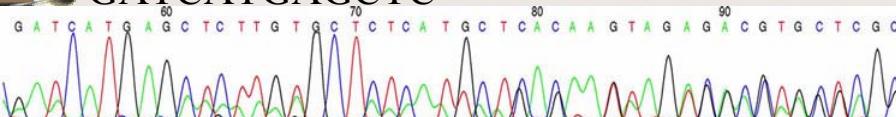
-
- Phil Green
 - Genome Sciences, University of Washington
 - <http://www.phrap.org>
 - Provides software and documentation

Phred

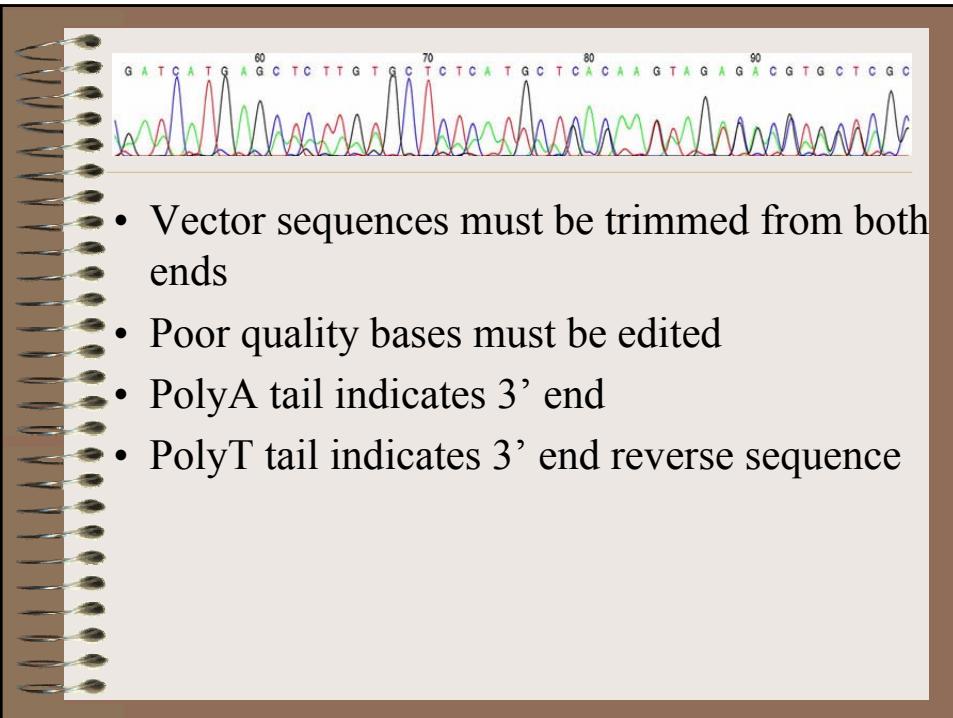
- Software reads sequencing trace files
- Calls bases
- Assigns a quality value to each called base
 - Correct and incorrect base calls
 - Quality values allow sequence trimming
- Works with Amersham Biosciences, Applied Biosystems, Beckman, LI-COR Life Sciences instruments

Which base reads are reliable

• GATCATGAGCTC



Phred



5' end
 TTTATCATGGCTGCCCTAGGGCGAT
 GAATGATCGTATGCCAGCTAAAAAAA
 AAAATCCGCCG
 3' end

From 5' end:
 ATG = methionine - possible start site
 TGA=STOP site
 AAAA...= possible polyA tail
 Remaining 3' sequence may be cloning vector sequence

Phrap

- Assembling shotgun DNA sequence data
- Improves assembly accuracy in presence of repeats
- Provides extensive assembly information to assist in trouble-shooting assembly problems
- Handles large data sets

Consed

- Automatically chooses finishing reads
- Speeds up finishing
- Integrated with Phrap
- DNA editing more efficient

The screenshot shows a Microsoft Internet Explorer window with the title "SWBIC - DNA Sequencing Software - Microsoft Internet Explorer". The address bar contains the URL <http://www.swbic.org-links/1.4.1.php>. The page content is titled "DNA Sequencing Software" and includes a navigation menu with tabs for "Educational Resources", "Internet Resources", and "Products & Services". On the left, there is a sidebar with links to "iDNAfication", "DoD Biotech Resources", "Minority Student Resources", "Bioinformatics Tools", and "Search Internet Resources". The main content area displays a list of software programs:

- BASS**: [Whitehead Institute for Biomedical Research] A software program for tracking, extracting, and base calling DNA sequencing gels.
- Chromas1.2**: [Technelysium] This program displays and prints chromatogram files from ABI automated DNA sequencers and Staden SCF files, and it allows the user to manipulate the sequences.
- DNATools**: [Carlsberg Laboratory] Software for handling and analysis of nucleotide and protein sequences. This downloadable package also has special functions for EST and SAGE

At the bottom of the page, there is a link "[more info]" next to each software entry.

The screenshot shows a web page with the heading "Sequence from your cDNA clone". Below the heading is a horizontal line. The sequence itself is a long string of DNA bases:

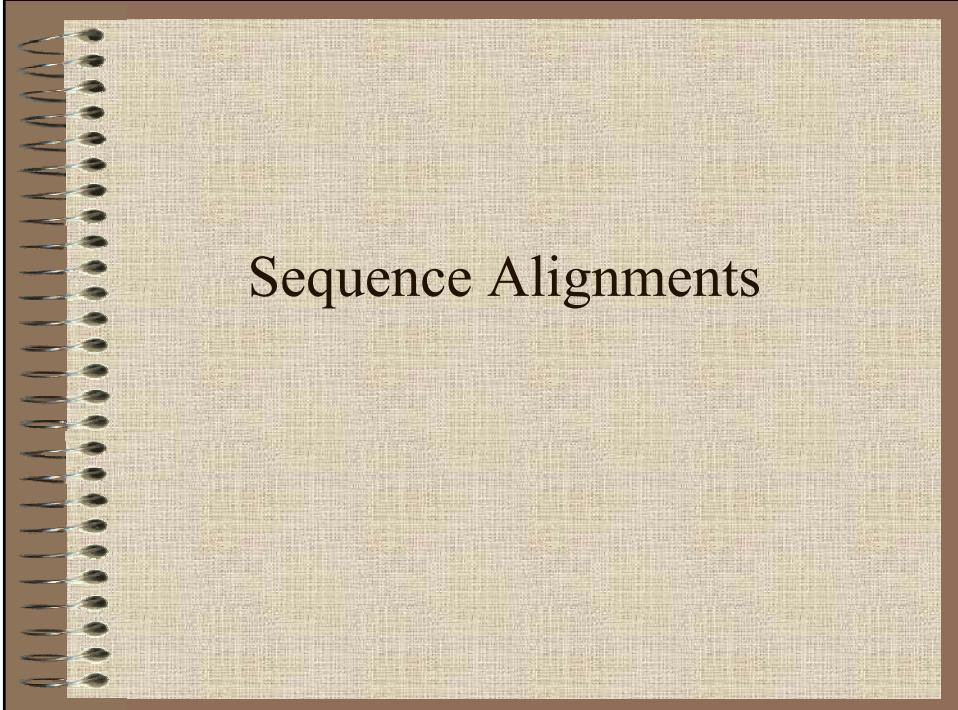
```
TACAGGGGTCCCCCTCCGGCGTCGGCTTCTCGATTCCAAGGGGAATGTTTT  
AAAGGCTTACATTGAGTCCGCTGCTTATAACCCAGCTGGGACCGCTTCA  
GGCCGCCATCGTCGCCCTCATGCCGGCGGGTGGGATTATGAAGAGATTG  
TTGCGGCGGTGTTGGAGAAGGAAGGGGCGGTCAAAACAGGATCACAC  
TGCAAGGTTGCTGCTCATTCCATAGCGCCACGCTGCCACTAACAAATTCT  
TGCTTCTCAATCTC
```

DNA sequence assembly

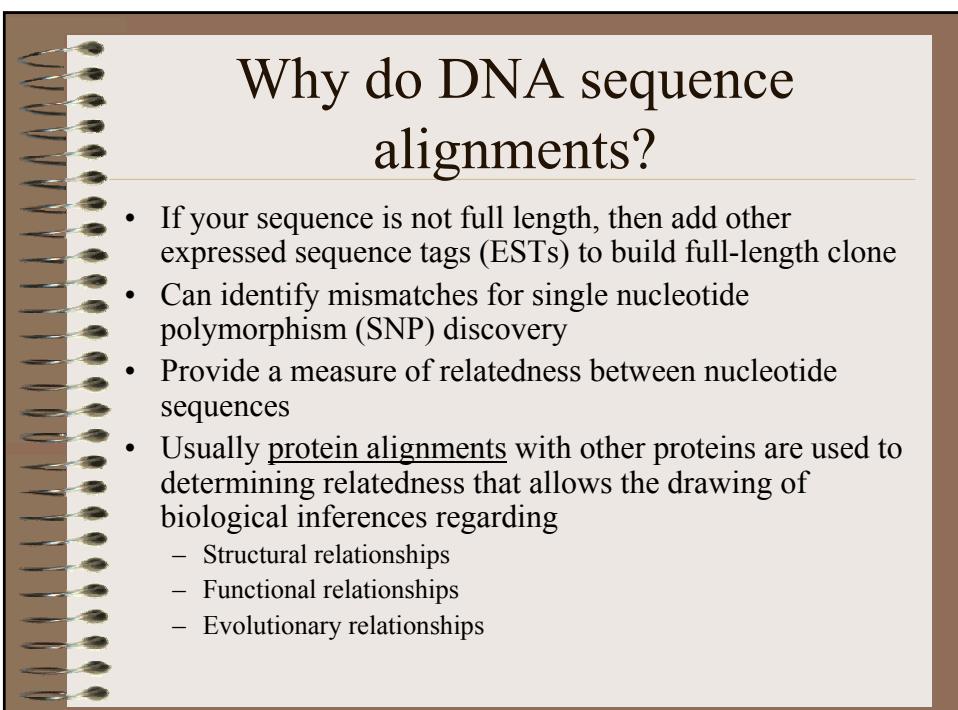
- PHRAP
- ConSED
- CAP3
- DNASTAR by Lasergene
 - Commercial - <http://www.dnastar.com/>
- Sequencher- commercial automated sequencers
 - <http://www.genecodes.com/>
 - Sequencher protocol
 - <http://bip.weizmann.ac.il/sequencher/sequencher.html>

cystein protease

```
GGAGCTCCACCGCGGTGGCGGCCGTTCTAGAACTAGTGGATCCCCGGCTGCAGGAATTGGCACAGAACAGTGG  
GAG  
GGAGATCCA AAAAAGAGAGATGGAAAAGATGGCGCGTGTGATCACAGTGGTGTGGCGGTGGGGGTGCTATTATG  
CG  
CCGGGGCGGTGCGTGTGGTGGAGGGCGCGAACCCGATACGAATGGTGTGGCGTGGAGGGGAGGTGGITC  
GG  
GTGATCGGGGAGTGCCGGCGTGTGAAGTTGCTAGGTTCTGAGCAGGTTGGGAAGAGTTACCAAAGCGAGGAA  
GA  
GATGAAGAGAGGTACGGAGATACTCGAGAACATCTCAGGTTATCCGCTCCCACAACAAGAGCTTGGCCCTATACTC  
T  
CTCTGTTAATCATTTGCTGATGGACTTGGAGGATTCAAAAGACACAGACTAGGGAGCTGCCAAAATGCTCTGCC  
A  
CTCTTAACGGCAACCACAAGCTACCGATGCTGTTCTCTCCAAACGAAAGACTGGAGAAAAGAGGTATAGTGAGTT  
CA  
GTTAAAGATCAAGGCAGCTCGGGATCATGCTGGACATTCAAGACAACGGGGCTTAAAGCAGCCTATGCAAGCA  
TT  
TGGGAAGAGTATCTCTTCTGAGCAGCAGTAGTGGACTGTGCTGGCCCTTCAACAACTTGGCTGCCATGGTGGG  
T  
TGGCCATCACAAGCCTTGTAGTACATTTAAACATGGTGGACTAGAGACAGAGGAAGCATATCCCTACACAGGAAAAG  
AT  
GGTGTGCAAATTCAGCTGAAAATGGTGTGTTCAAGTCAGCTTGTGAAATATCACCTGGGTGCTGAAGATG  
A  
ACTAAACATGCAGTGCATTGTCGGCCAGTTAGTGTGGCCCTTCAAGTGGTGAATGGGTTCAATTTCAGGAAAT  
G  
GAGTTTCACTAGTGACACTGTGGTAGCAGTCCCAAGGATGTAACCATGCCCTTGTGCTGTTGGATATGGAGTTGA  
A  
AATGGTGTCCCATATTGGCTATAAAAAAATCATGGGGAGAAAAGCTGGGGAAAATGGCTACTCAAGATGGAATTG  
GG  
GAAGAACATGTGGTGTGCAACTTGTGCACTTATCCAATTGGCATAAAATTGCAAAATATGGCCCTGGTGA  
C  
TACCACTTGTGTCAGAGTTAGAGCTATTGCTGATGCCAGTATGATGAAATGATGATGATTTAAAGATAAGTGTAA  
T  
TGATGATGAAAATTGCTCTAGTGGGTGGCATGATGTTAAAAGCTAGAATGTTGTAATACACATAAGTAT  
A  
TTATGGCTTAAATGTTGTTACAGACATAAAACGATCATATTGATAGTCAAGTTACATATTGTTATTGATTG  
ATGCTCCGCTTAAATACAGTTAAAGATGAGCTTGTGCTACTTGTGCACTATGCAACACATT
```



Sequence Alignments



Why do DNA sequence alignments?

- If your sequence is not full length, then add other expressed sequence tags (ESTs) to build full-length clone
- Can identify mismatches for single nucleotide polymorphism (SNP) discovery
- Provide a measure of relatedness between nucleotide sequences
- Usually protein alignments with other proteins are used to determining relatedness that allows the drawing of biological inferences regarding
 - Structural relationships
 - Functional relationships
 - Evolutionary relationships

Similarity

- A quantitative measure
- Based on an observable
- Usually expressed as percent identity
- Quantifies changes that occur as two sequences diverge
 - Substitutions
 - Insertions
 - Deletions
- Identifies residues crucial for maintaining a protein's structure or function

Similarity

- High degrees of similarity *might* imply
 - A common evolutionary history
 - A possible commonality in biological function



Homology

- Implies an evolutionary relationship
- May apply to the relationship
 - Between genes separated by the event of speciation (orthology), ie. orthologous genes
 - Between genes separated by the event of genetic duplication (paralogy), ie. paralogous genes

- 
- Orthologs
 - Sequences are direct descendants of a sequence in a common ancestor
 - Most likely have similar domain structure, three dimensional structure, and biological function
 - Paralogs
 - Related through a gene duplication event
 - Provides insight into evolution, ie. adapting a pre-existing gene product for a new function



Global Sequence Alignments

- Sequence comparison along the entire length of two sequences being aligned
- Best for highly similar sequences of similar length
- As the degree of sequence similarity declines, global alignment methods tend to miss relationships

Local Sequence Alignments

- Sequence comparison intended to find the most similar regions in two sequences being aligned
- Regions outside the area of local alignment are excluded
- More than one local alignment could be generated
- Best for sequences that share some similarity or for sequences of different lengths

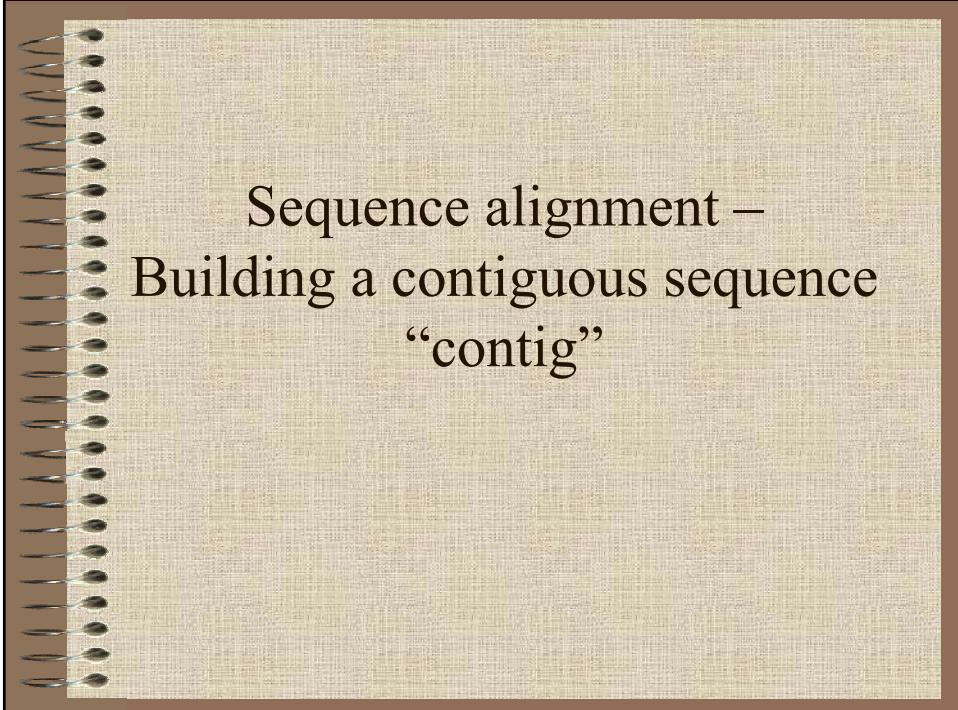
Scoring Matrices

- Empirical weighting scheme to represent biology
- DNA only has A,T,G,C
- Protein has amino acids; relatedness among amino acids; function; charges; side groups

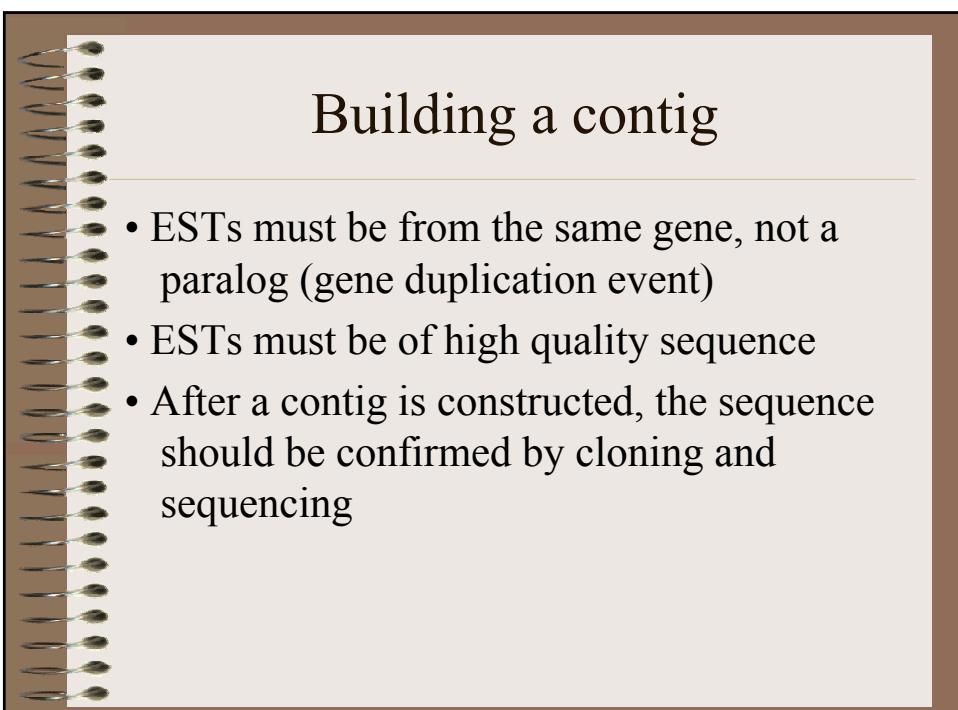
Matrix Structure: Nucleotides

A	T	G	C	A	T	G	C	A	T	G	C	A	T	G	C	A	T	G	C
A	-4	-4	-4	-4	1	1	-4	-4	1	1	-4	-1	-1	-1	-1	-4	-1	-1	-1
T	-4	1	-4	-4	1	-4	1	1	1	-4	-1	-1	-1	-1	-1	-4	-1	-1	-1
G	-4	1	-4	-4	1	-4	1	-4	1	1	-4	-1	-1	-1	-1	-4	-1	-1	-1
C	-4	-4	-4	1	1	-4	1	-4	1	1	-4	-1	-1	-1	-1	-4	-1	-1	-1
A	-4	-4	1	1	-1	-4	-2	-2	-2	-2	-1	-1	-1	-1	-1	-3	-3	-3	-1
T	1	1	-4	-4	-1	-2	-2	-2	-2	-2	-3	-3	-3	-3	-3	-1	-1	-1	-1
G	1	-4	1	-4	-2	-2	-1	-4	-2	-2	-2	-3	-2	-2	-2	-3	-1	-1	-1
C	-4	1	-4	1	-2	-2	-1	-2	-1	-2	-2	-3	-1	-1	-1	-3	-1	-1	-1
A	-4	1	-4	1	-2	-2	-1	-2	-1	-2	-2	-3	-1	-1	-1	-3	-1	-1	-1
T	-4	1	-4	1	-2	-2	-1	-2	-1	-2	-2	-3	-1	-1	-1	-3	-1	-1	-1
G	1	-4	1	-4	-2	-2	-1	-2	-1	-2	-2	-3	-1	-1	-1	-3	-1	-1	-1
C	-4	-1	-1	-1	-1	-2	-1	-1	-1	-1	-2	-3	-1	-1	-1	-3	-1	-1	-1
A	-1	-4	-1	-1	-1	-2	-1	-1	-1	-1	-2	-3	-1	-1	-1	-3	-1	-1	-1
T	-1	-1	-4	-1	-1	-2	-1	-1	-1	-1	-2	-3	-1	-1	-1	-3	-1	-1	-1
G	-1	-1	-1	-4	-1	-1	-1	-1	-1	-1	-2	-3	-1	-1	-1	-3	-1	-1	-1
C	-2	-2	-2	-2	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1

- Simple match/mismatch scoring scheme
- Assumes each nucleotide occurs 25% of the time



Sequence alignment – Building a contiguous sequence “contig”



Building a contig

- ESTs must be from the same gene, not a paralog (gene duplication event)
- ESTs must be of high quality sequence
- After a contig is constructed, the sequence should be confirmed by cloning and sequencing



EST 1: gaggctatgccgtccgagattacggcgttacaggattcagatt
EST 2: ggacccaagttcacgtccaatattgtgttgaccatagaaaaaaaaaa
EST 3: acggcgttacaggattcagattcatggacccaagttcacgtc

EST alignment to make contig

EST 1: gagectatgecgtecgagattacggcgttacaggattcagatt

EST 3: acggcgttacaggattcagattcatggacccaagttcacgtc

EST2: ggacccaagttcacgtccaatattgtgttgacc
atagaaaaaaaaaa

Consensus sequence:

gagectatgcgtccgagattacggcgttacaggattcagattcatggacccaagttcacgtccaatattgtgttgaccata
aaaaaaaaaa

In this example EST 3 forms a bridge to connect EST 1 and EST 2

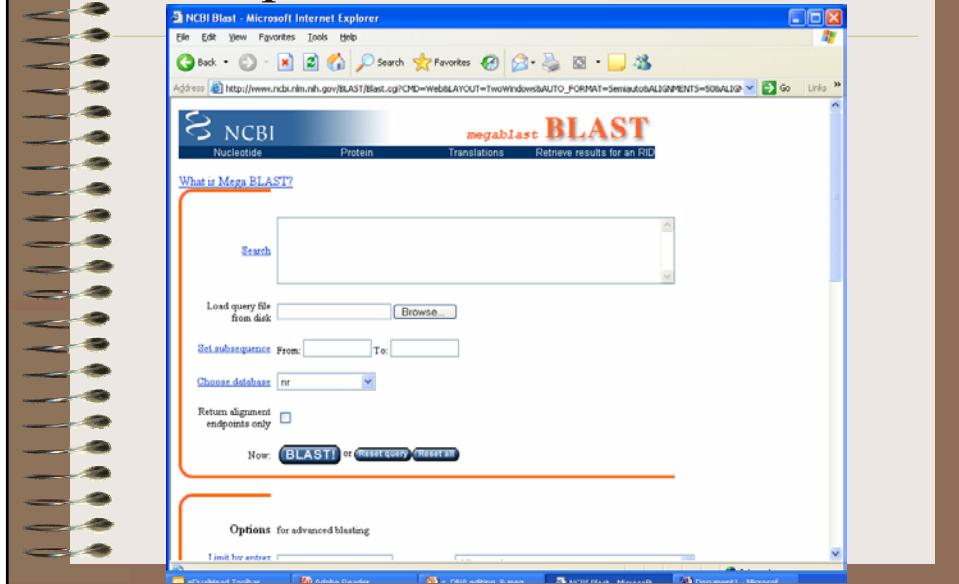
DNA STAR EXAMPLE

Making a contig from EST sequences

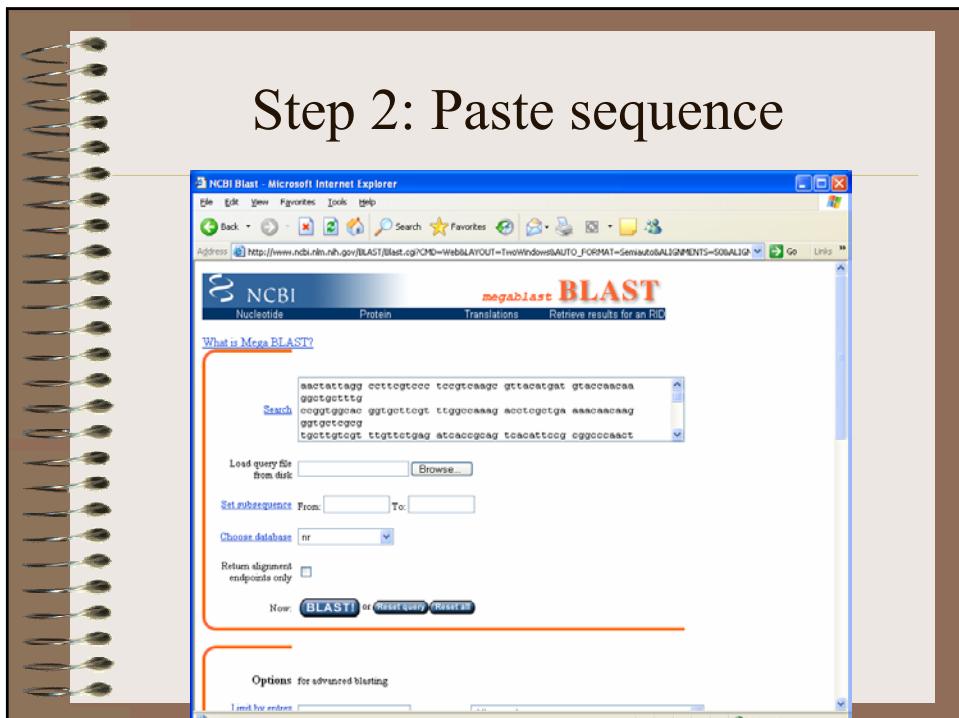
This is your sequence from a clone

```
1 aactattagg ccttcgtccc tccgtcaagc gttagatgtt gttaccaacaa ggctgccttg  
61 ccgggtggcac ggtgcgttgcgt ttggccaaag acctcgctga aaacaacaag ggtgctcgcg  
121 tgcttgtcgat ttgttctgat atcaccgcag tcacatccg cggcccaact gacacccatc  
181 ttgatagcct tggggtaaa gccttgcgtt gatgtttgc agccgcgttc attgttggat  
241 cagaccctt accagttgaa aaggcttgcgtt tcacgttat ctggactgccc caaacaatcc  
301 ttccagacag tgaaggggctt attgtatggcc accttcgcga agttggactc actttccatc  
361 tcctcaagga tttcccttgcata cttcatctca agaatattga gaaggccctt gttgaaggct  
421 tccaaccctt gggaaatctcc gattacaatt ctatcttcgtt gattgcacac cct
```

Step 1: Go to BLAST search



Step 2: Paste sequence



Step 3: Set constraints and options

actttccatc
ttcctcaagg tgcgttggaa
ctcatcttca agatatttga
gaaaggccgtt
gttggacgact
tcacaccgtt
gggaaatctcc
gattacaatt
cttatcttctg
gatttgacac
cct

Search

Set subsequence From: [] To: []

Choose database nr

Now: **BLAST!** or [Reset query](#) [Reset all](#)

Options for advanced blasting

Limit by entrez query or select from: All organisms

Choose filter Low complexity Human repeats Mask for lookup table only Mask lower case

Expect 10

Load query file from disk [Browse...](#)

Set subsequence From: [] To: []

Choose database test

Return alignment endpoints only

Now: **BLAST!** or [Reset query](#) [Reset all](#)

Options for advanced blasting

Limit by entrez query or select from: All organisms

Choose filter Low complexity Human repeats Mask for lookup table only Mask lower case

Expect 10

Word Size 28

Percent Identity match_mismatch scores None, 1,-2

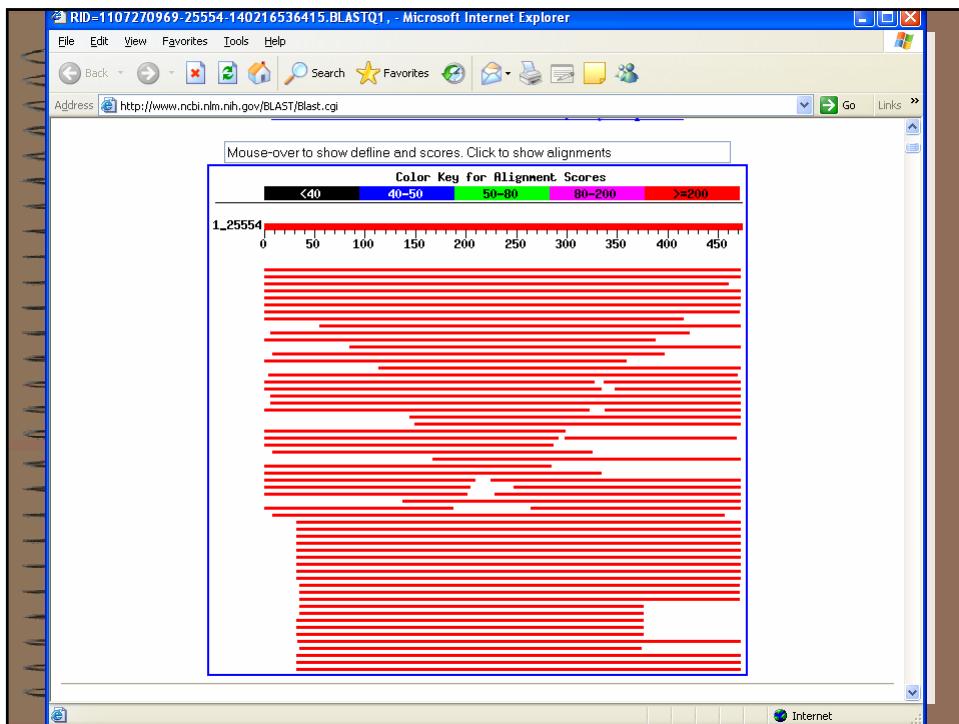
Step 4: BLAST!

The screenshot shows the NCBI BLAST interface in Microsoft Internet Explorer. The title bar reads "NCBI Blast - Microsoft Internet Explorer". The address bar shows the URL http://www.ncbi.nlm.nih.gov/BLAST/Blast.cgi?CMD=Web&LAYOUT=TwoWindows&AUTO_FORMAT=SemiAuto&ALIGNMENTS=50&ALIQ=1. The main content area is titled "Format". It contains various configuration options:

- Show: Graphical Overview Linkout Sequence Retrieval NCBI-g
- Alignment:
- Use new formatter: Masking Character Default (X for protein, n for nucleotide) Masking Color Black
- Number of Descriptions: 100 Alignments: 50
- Alignment view: Hit Table
- Start formatting from query #:
- Limit results by Entrez query: [] or select from: All organisms
- Expect value: [] []
- Layout: Two Windows
- Autoformat: Semi-auto
- Results file: []

At the bottom left is a blue button labeled "BLAST!" with a red arrow pointing to it. Below the button is the text "Get the URL with preset values?

The screenshot shows the NCBI BLAST interface in Microsoft Internet Explorer. The title bar reads "NCBI Blast - Microsoft Internet Explorer". The address bar shows the URL <http://www.ncbi.nlm.nih.gov/BLAST/Blast.cgi>. The main content area is titled "formatting BLAST". It displays the message: "Your request has been successfully submitted and put into the Blast Queue." Below this, it says "Query = (473 letters)" and "Your search was limited by an Entrez query: Glycine max". A text input field shows "The request ID is: 1107270343-22553-173958027945 BLAST01". Below the input field is a blue button labeled "Format" with a red arrow pointing to it. To the right of the button is the text "The results are estimated to be ready in 1 minutes 20 seconds but may be done sooner". Further down, it says "Please press 'FORMAT!' when you wish to check your results. You may change the formatting options for your results via the form below and press 'FORMAT!' again. You may also request results of a different search by entering any other valid request ID to see other recent jobs." At the bottom is a "Format" configuration panel identical to the one in the first screenshot.



RID=1107270969-25554-140216536415.BLASTQ1, - Microsoft Internet Explorer

File Edit View Favorites Tools Help

Address http://www.ncbi.nlm.nih.gov/BLAST/Blast.cgi

Sequences producing significant alignments:

	Score	E
	(bits)	Value
gi 13480079 gb BG509422.1	938	0.0
gi 37996764 gb CF808353.1	922	0.0
gi 20813264 gb BQ297742.1	906	0.0
gi 37996212 gb CF807801.1	882	0.0
gi 15815451 gb BT787726.1	882	0.0
gi 37996627 gb CF808216.1	866	0.0
gi 23728449 gb BU762277.1	864	0.0
gi 37995445 gb CF807034.1	793	0.0
gi 37995974 gb CF807563.1	765	0.0
gi 15337571 gb BT498227.1	763	0.0
gi 33390507 gb CA853702.1	729	0.0
gi 16346573 gb BI1972168.1	722	0.0
gi 15813558 gb BT785833.1	668	0.0
gi 37994482 gb CF806228.1	664	0.0
gi 27427571 gb CA939091.1	664	0.0
gi 19936478 gb BQ080882.1	656	0.0
gi 19936194 gb BQ080763.1	618	e-176
gi 37994162 gb CF805908.1	609	e-173
gi 19934725 gb BQ079755.1	609	e-173
gi 8402207 gb BE057841.1	609	e-173
gi 19938183 gb BQ081600.1	601	e-171
gi 13478388 gb BQ507884.1	599	e-170
gi 15287478 gb BI1471369.1	595	e-169
gi 17400989 gb BM177771.1	569	e-161

NCBI Sequence Viewer v2.0 - Microsoft Internet Explorer

File Edit View Favorites Tools Help

Back Stop Home Search Favorites Mail Print Copy Paste Links

Address http://www.ncbi.nlm.nih.gov/entrez/query.fcgi?cmd=Retrieve&db=Nucleotide&list_uids=13480079&dopt=GenBank Go Links

Nucleotide

PubMed Nucleotide Protein Genome Structure PMC Taxonomy OMIM Books

Search: Nucleotide for Go Clear

Limits Preview/Index History Clipboard Details

Display GenBank Send all to file

Range: from begin to end Reverse complemented strand Features: SNP CDD MGC HPRD

1: BG509422 Reports sad13f03.y1 Gm-c1074 [gi:13480079] Links

LOCUS BG509422 473 bp mRNA linear EST 24-JUL-2004

DEFINITION sad13f03.y1 Gm-c1074 Glycine max cDNA clone GENOME SYSTEMS CLONE

ID: Gm-c1074-246 5' similar to SW:CHS1_SOYBEAN P24826 CHALCONE SYNTHASE 1 ; mRNA sequence.

ACCESSION BG509422

VERSION BG509422.1 GI:13480079

KEYWORDS EST.

SOURCE Glycine max (soybean)

ORGANISM Glycine max
Eukaryota; Viridiplantae; Streptophyta; Embryophyta; Tracheophyta; Spermatophyta; Magnoliophyta; eudicotyledons; core eudicots; rosids; eurosids I; Fabales; Fabaceae; Papilionoideae; Phaseoleae; Glycine.

REFERENCE 1 (bases 1 to 473)

AUTHORS Shoemaker,R., Keim,P., Vodkin,L., Erpelding,J., Coryell,V., Khanna,A., Boilla,B., Marra,M., Hillier,L., Kucaba,T., Martin,J., Beck,C., Wylie,T., Underwood,K., Steptoe,M., Theising,B., Allen,M., Bowers,Y., Person,B., Swaller,T., Gibbons,M., Pape,D., Harvey,N., Schurk,R., Ritter,E., Kohn,S., Shin,T., Jackson,Y., Cardenas,M., McCann,R., Waterston,R. and Wilson,R.

TITLE Public Soybean EST Project

JOURNAL Unpublished (1999)

NCBI Sequence Viewer v2.0 - Microsoft Internet Explorer

File Edit View Favorites Tools Help

Back Stop Home Search Favorites Mail Print Copy Paste Links

Address http://www.ncbi.nlm.nih.gov/entrez/query.fcgi?cmd=Retrieve&db=Nucleotide&list_uids=13480079&dopt=GenBank Go Links

tissue with *Pseudomonas syringae* pv. glycinea carrying the avrB gene (Genetics 141:1597-1604). Plant tissue (expanded unifoliate leaves) was collected at 2, 4, 8, 12, 24, 36, and 53 hrs after inoculation and their mRNA pooled equally for cDNA construction. The library was prepared using the Stratagene pBluescript II SK(+) library construction kit. Complementary DNA was synthesized from mRNA using a primer consisting of a poly(dT) sequence with an XbaI restriction site. EcoRI adaptors were ligated to the blunt-ended cDNA fragments followed by XbaI digestion. The cDNA insert is protected from XbaI digestion via methylation during first strand synthesis. The cDNA fragments were directionally cloned into the EcoRI-XbaI restriction site of the pBluescript vector. The ligated cDNA fragments were transformed into E.coli ElectroMax DH10B host cells. Plant care, inoculations, and library construction were performed by Steve Clough (Lila Vodkin lab, University of Illinois)."

ORIGIN

```
1 aactatttagg cttctgtccc tccgtcaacg gttatcatgt gtaccaacaa ggctgttgg
 61 ccgggtggcac ggtgttgtcg ttggccaaag acctcgctga aaacacaacg gggtgttcgcg
 121 tggtgttgct ttgttgtcgat atcaccccgat tcataatccgg cggcccaact gagaccatcc
 181 ttgatagcct ttgtggtaaa gccttggttt gagatgtgtc agccgcgttc attgttgtggat
 241 cagaccctt accaggtaaa aagccgtttt ttcaatgtttt ctggactgcc caaacaatcc
 301 ttccagacag tgaaggggctt atttatgtggcc acctttcgca agttggactt acttttcatc
 361 ttctcaaggg tggttccgttcca ctcatcttca aagatattga gaaggcttgg gttgaaggctt
 421 tccaaaccctt gggaaatctcc gattacaattt ctatcttctt gattgcacac cct
//
```

[Disclaimer](#) | [Write to the Help Desk](#)
[NCBI](#) | [NLM](#) | [NIH](#)

Jan 27 2005 17:14:21

Collect sequences into a series of files so they can be aligned

File 1.

```
1 aactattagg ccttcgtccc tccgtcaagc gtacatcatgt gtaccaacaa ggctgcttg  
61 ccgggtggcac ggtgttcgt ttggccaaag acctcgctga aaacaacaag ggtgtcg  
121 tgcttgcgt ttgttcgtag atcacccgcag tcacattccg cggcccaact gacaccatc  
181 ttgtatggct tgggggtcaa gcctgtttg gagatgggtc agccgcgtc atttgtggat  
241 cagacccctt accagtgaa aagcccttgc ttcaagcttat ctggactgcc caaacaatcc  
301 ttccagacag tgaaggggctt attgtggcc accttcgcga agtgtggactc actttccatc  
361 tcccaagga tggttgcgtt ctcatctca agaatattga gaaggccctt gttgaaggct  
421 tccaaccctt gggaaatctcc gattacaatt ctatctctg gattgcacac cct
```

File 2

```
1 ggcaatcaag gaatggggtc aacccaagtc caagattacc catctcatct ttgcaccac  
61 tagtgggtc gacatccgtt gtgtgttca tcagtcactt aaactattag gcccgttc  
121 ctccgtcaag cggttcaatgtt tgatccaaaca aggctgtttt gcccgtggca cgggtgc  
181 ttggccaaa gacccgtctg aaaacaacaa ggggtcgcc gtgcgtcg ttgttctga  
241 gtacccgcga gtccatcc gggcccaat tgacacccat ttgtatggcc ttgtgggtca  
301 agcctgtttt ggagatggtg cagccgtgtt cattgttggta tcagacccat taccagtgt  
361 aaaggcccttgg ttcagctta tctggactgc ccaaaaatcc tcctcagaca gtgaaggggc  
421 tattgtatggc cacccttcgcga aagttggact cactttcat tcctcaagg atgttccctgg  
481 actcatcttca aagaataatggc aagaaggctt ggttggcc ttccaaacctt tggttgc  
541 cgattacaat tctatctctt ggattgcaca ccctgggttgc cccgcaattt tggttgc  
601 tgaggcttgg ttaggttgcgtt ggcctgaaa aatggaaatggctt actagacatg tgctc  
661 gtatgttacatg
```

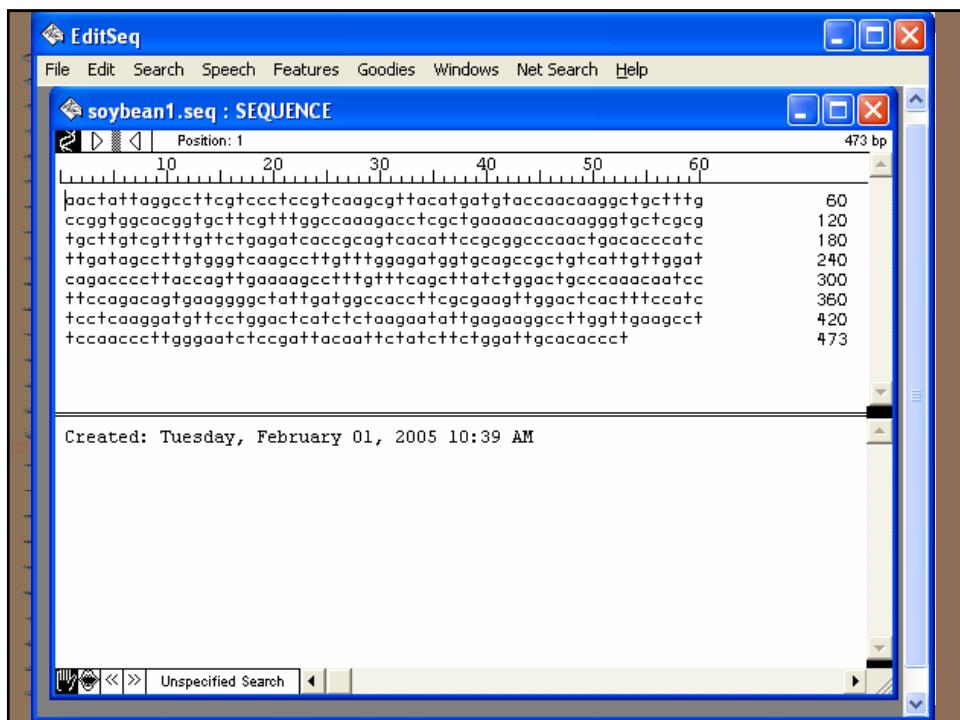
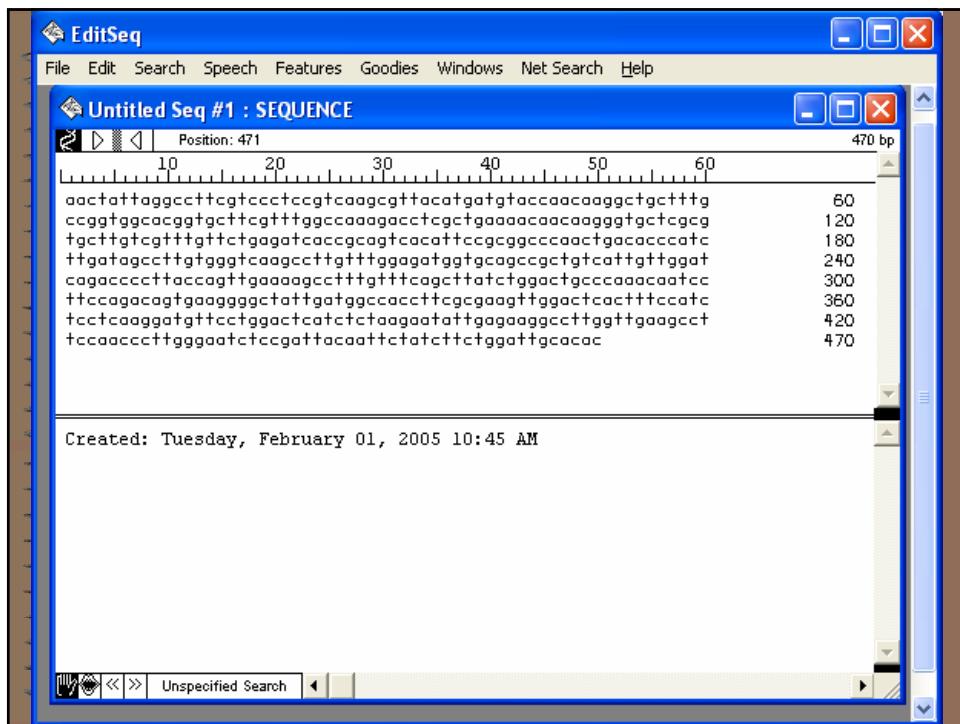
File 3

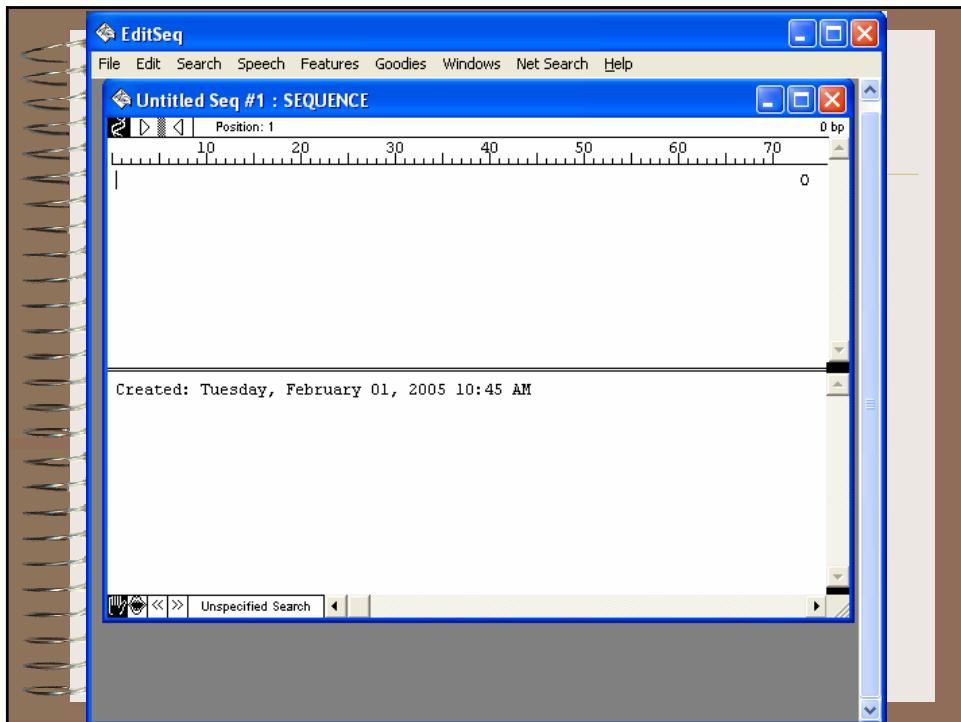
```
1 ttttttttttttttttttttttttttttttttttttttttttttttttttttttttttttttttttttttttt  
61 accccaagtcc aagatttaccctt atctcatctt ttgcaccactt agtgggtgc acatgc  
121 tgctgattttt cagcttcaactt aacttattttt ctttcgtccc tccgtcaagc ttgtatgtat  
181 gtaccaacaa gggtgtttt ccgggtggcac ggtgtcgcc ttggccaaag acctcg  
241 aaacaacaag ggtgtcgcc tgctgtcgat ttgttgcgtt atcacccgcag tcacattcc  
301 cggcccaactt gacacccatc ttgtatggctt ttgggttcaaa gcccgtttt gggatgggt  
361 agccgcgttc atttgtggat cagacccctt accagtgttggaa aagcccttgc ttcaagcttat  
421 ctggacttgc cccatccatcc ttccatccatc gtttggggctt attgtggcc accttcgc  
481 agttggactc actttccatc ttccatccatc gtttggggctt attgtggcc accttcgc  
541 gaaggcccttgg ttaggttgcgtt ggcctgaaa aatggaaatggctt actagacatg tgctc  
601 g
```

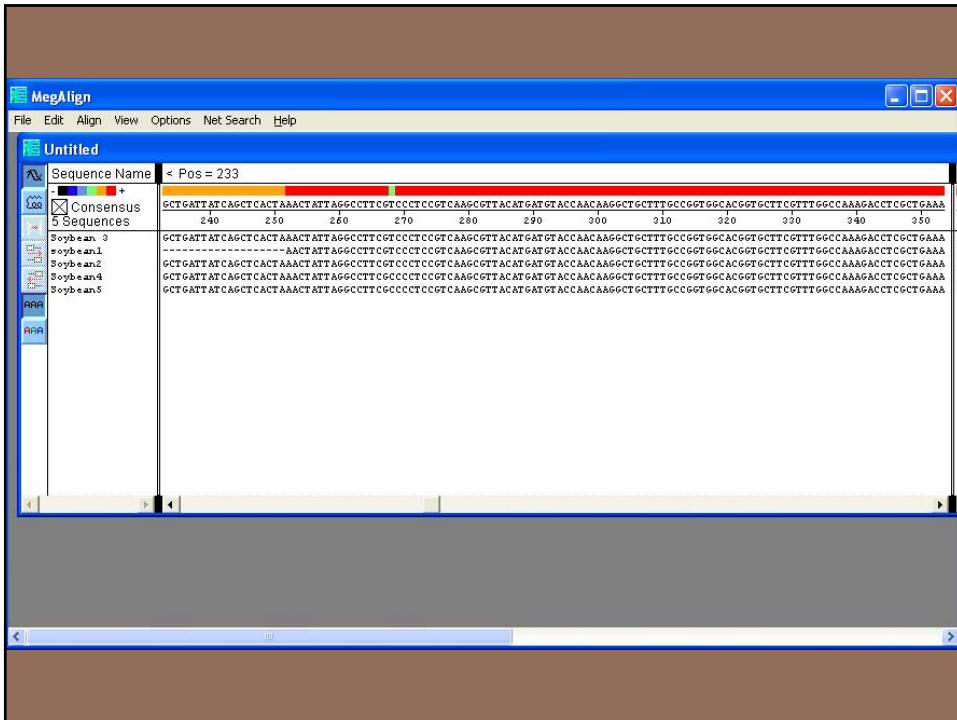
DNA STAR

- Edit sequence
 - Allows you to import and edit DNA and protein sequences
- Megalign
 - Allows you to align DNA and protein sequences

Edit Sequence







CAP EST Assembler

Contig Assembly Program

- <http://bio.ifom-firc.it/ASSEMBLY/assemble.html>
- Can use up to 30,000 EST sequences
- Fragment maximum is 30,000 bp
- Sequences must be in FASTA format
- Huang, X. 1992. Genomics 14: 18-25
- Huang, X. 1996. Genomics 33: 21-31

Edit file

- Some DNA and protein alignment software requires a specific format
- FASTA format
 - Header HAS TO start with ‘ > ’
 - A description should follow
 - For DNA only five letters A,C,T,G,N allowed
 - No numbers

```
>soybean1
ATTCCCTTAGGATC...
>soybean 2
TCCGTCAAGGTGTT...
>soybean3
GGCTATGGCCTAAT...
```

The screenshot shows a Microsoft Internet Explorer window with the title "The CAP Sequence Assembly Machine - Microsoft Internet Explorer". The address bar displays the URL <http://bio.ifom-firc.it/ASSEMBLY/assemble.html>. The main content area is titled "The CAP EST Assembler at IFOM". Below the title, there is a bulleted list of assembly parameters:

- Maximum sequence length for each sequence is 30.000 - Maximum number of sequences 10.000
- Timeout for interactive assembly is 5 minutes - Maximum uploadable data is 1 Megabyte
- MAXIMUM FRAGMENT LENGTH IS 30.000 bp - this software is optimized for EST assembly

Below the parameters, there is a text input field with the placeholder text "Enter sequences to assemble below, in [FASTA format](#):".

THE CAP Sequence Assembly Machine - Microsoft Internet Explorer

File Edit View Favorites Tools Help

Address: <http://bio.ifom-firc.it/ASSEMBLY/assemble.html>

Links >

- Maximum sequence length for each sequence is 30.000 - Maximum number of sequences 10.000
- Timeout for interactive assembly is 5 minutes - Maximum uploadable data is 1 Megabyte
- MAXIMUM FRAGMENT LENGTH IS 30.000 bp - this software is optimized for EST assembly

Enter sequences to assemble below, in **FASTA** format:

```
>Soybean1
aatggggtcacccaaagtccaaagattaccatctatctttgcaccactagtgggtc
gacatgcctgtgtgttgcattat
cagtcactaaactattaggccctccgcctcgtaagcggttacatgtatgtaccaaca
aggctgccttgcgggtggcac
gggtcttcgtttggccaaagacccgtgtgaaaacaacaagggtgtcgctgtcg
tttgttcttgatgtacccgcag
tcacatttcggccaaactgtacacccatcttgatagccttgggtcaagccttgc
ggagatgttgcacgcgcgtc
atgttggatgtacaccccttaccgttggaaaagccttgcgtttagcttgcactgc
ccagacaatcttcagacag
tgaaggggctattttgtggacaccccttcgcgaagggttctcatctcccaagg
atgttctggactatctcca
agaatattgttggaccccttgggttgaagcccttcaacccttggaaatctccgattacaat
tctatcttcggatttgcacac
cct
>Soybean2
ggcacatcaaggatgggtcaacccaaagtccaaagattaccatctatctttgcacca
ctagtgtgtgcacatgtcgt
gtgtgttattatcgtactaaactattaggccttgcgtccctccgtcaagcgttacatg
```

Internet

Caps Assembler at IFOM Output - Microsoft Internet Explorer

File Edit View Favorites Tools Help

Address: <http://bio.ifom-firc.it/cgi-bin/Assembly/capassembly.pl>

The Cap Sequence Assembler at IFOM

Program Cap3 used

Assembly results:

Number of segment pairs = 20; number of pairwise comparisons = 10
 '+' means given segment; '-' means reverse complement

Overlaps	Containments	No. of Constraints	Supporting Overlap
***** Contig 1 *****			
SOYBEAN4+	SOYBEAN1+ is in SOYBEAN4+ SOYBEAN5+ is in SOYBEAN1+ SOYBEAN3+ is in SOYBEAN4+		
SOYBEAN2+			
DETAILED DISPLAY OF CONTIGS	***** Contig 1 *****		
SOYBEAN4+	GTCGATGATTAAGAACGGATACATGTACTTAAACGAAGAGATCCTGAAAGAGAATCCGAG		
consensus	GTCGATGATTAAGAACGGATACATGTACTTAAACGAAGAGATCCTGAAAGAGAATCCGAG		
SOYBEAN4+ SOYBEAN3+	TGTTTGCTTACATGGCACCTCGTTGGATGCAAGGCAAGACATGGTGGTTGTGGAGGT TGTGGAGGT		
consensus	TGTTTGCTTACATGGCACCTCGTTGGATGCAAGGCAAGACATGGTGGTTGTGGAGGT		

Internet

Caps Assembler at IFORM OUTPUT - Microsoft Internet Explorer

File Edit View Favorites Tools Help

Address: http://bio.ifom-firc.it/cgi-bin/Assembly/capassemble.pl

DETAILED DISPLAY OF CONTIGS

***** Contig 1 *****

SOYBEAN4+	GTCGATGATTAAGAAGCGATACATGTACTTAAACGAAGAGATCCTGAAAGAGAAATCCGAG
consensus	GTCGATGATTAAGAAGCGATACATGTACTTAAACGAAGAGATCCTGAAAGAGAAATCCGAG
SOYBEAN4+	TGTTTGCTTACATGGCACCTCTGTTGGATGCAAGGCAAGACATGGTGGTTGTGGAGGT
SOYBEAN3+	TGTGGAGGT
consensus	TGTTTGCTTACATGGCACCTCTGTTGGATGCAAGGCAAGACATGGTGGTTGTGGAGGT
SOYBEAN4+	ACCAAAGTTGGAAAAAGAGGCTGCAACTAACCGCAATCAAGGAATGGGTCAACCCAAGTC
SOYBEAN1+	AATGGGGTCAACCCAAGTC
SOYBEAN5+	AATGGGGTCAACCCAAGTC
SOYBEAN3+	ACCAAAGTTGGAAAAAGAGGCTGCAACTAACCGCAATCAAGGAATGGGTCAACCCAAGTC
SOYBEAN2+	GGCAATCAAGGAATGGGTCAACCCAAGTC
consensus	ACCAAAGTTGGAAAAAGAGGCTGCAACTAACCGCAATCAAGGAATGGGTCAACCCAAGTC
SOYBEAN4+	CAAGATTACCCATCTCATTTTGCAACCACTAGTGGTGTGACATGCCCTGGTGTGATTA
SOYBEAN1+	CAAGATTACCCATCTCATTTTGCAACCACTAGTGGTGTGACATGCCCTGGTGTGATTA
SOYBEAN5+	CAAGATTACCCATCTCATTTTGCAACCACTAGTGGTGTGACATGCCCTGGTGTGATTA
SOYBEAN3+	CAAGATTACCCATCTCATTTTGCAACCACTAGTGGTGTGACATGCCCTGGTGTGATTA
SOYBEAN2+	CAAGATTACCCATCTCATTTTGCAACCACTAGTGGTGTGACATGCCCTGGTGTGATTA
consensus	CAAGATTACCCATCTCATTTTGCAACCACTAGTGGTGTGACATGCCCTGGTGTGATTA
SOYBEAN4+	TCAGCTCACTAAACTATTAGGCCTTCGCCCTCCGTCAAGCGTTACATGATGTACCAACA
SOYBEAN1+	TCAGCTCACTAAACTATTAGGCCTTCGCCCTCCGTCAAGCGTTACATGATGTACCAACA
SOYBEAN5+	TCAGCTCACTAAACTATTAGGCCTTCGCCCTCCGTCAAGCGTTACATGATGTACCAACA
SOYBEAN3+	TCAGCTCACTAAACTATTAGGCCTTCGCCCTCCGTCAAGCGTTACATGATGTACCAACA

Internet

What we learned today

- DNA editing
- Phred
- Phrap
- Consed
- DNA Sequencing software
- DNA sequence assembly
- Similarity searching with a DNA sequence
- BLAST